

Dr. (Amos) Xinshao Wang: Data + Modelling
Vision and Endeavours: To Design and Build
Data/Label/Time-Efficient, Robustly Reliable and Transparent
AI for Diverse Applications (Omics/Biosequences, Healthcare,
NLP, CV, etc).

Google Scholar: [yOBhB7UAAAAJ](https://scholar.google.com/citations?user=yOBhB7UAAAAJ)

Github: <https://github.com/XinshaoAmosWang>

Homepage+Blogs: <https://xinshaoamoswang.github.io/>

Email: xinshaowang@gmail.com

Phone: +44 (0) 7712 114316

2023-04-19 Wed, Invited Talk@Synteny

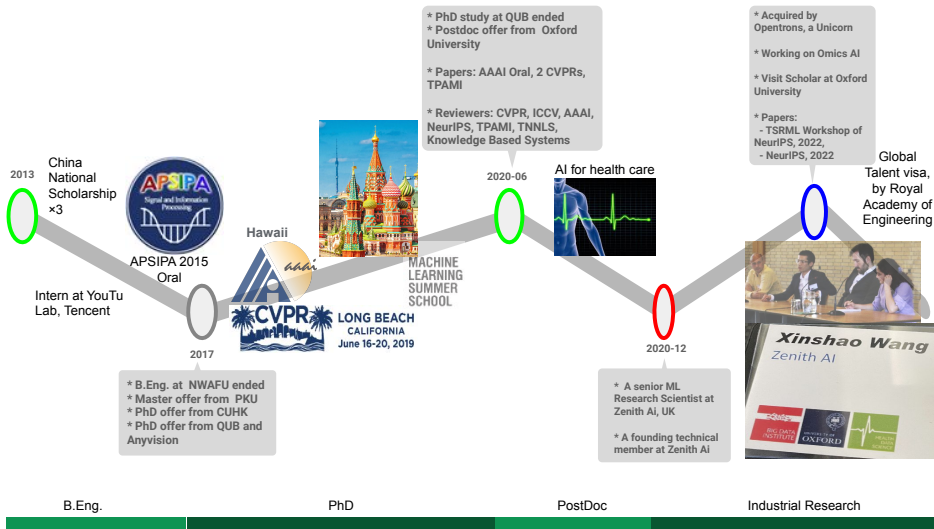
Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Self introduction: Overview



1 Self introduction

Overview

Work and experience

2 Technical novelty and principle: Example weighting is universal

3 Deep distance metric learning

Contribution: Structured learning on selected informative data

4 Robust deep learning (scope: transparent and interpretable ML)

Learning objectives

Contributions: an insightful finding and propose ProSelfLC

5 Omics AI

Cutting-edge OTS research: triplet design

Bioinformatics: alignment-based distance and virtual screening

Robust deep learning for protein modelling

The key: collecting, curating, and leveraging the data

6 Industrial R&D experience

Industrial R&D for real-world problems

R&D leading experience

Work and experience

- ① Deep distance metric learning
 - CVPR 2019 and TPAMI 2022
 - AAAI 2019 Oral
- ② Robust deep learning, model calibration and uncertainties:
 - CVPR 2021
 - Co-supervising a PhD student: [Trustworthy and Socially Responsible ML Workshop, NeurIPS 2022](#)
 - Trustworthy and Reliable ML Workshop 2023, ICLR “[IMAE for Noise-Robust Learning: Mean Absolute Error Does Not Treat Examples Equally and Gradient Magnitude’s Variance Matters](#)” (65 citations)
- ③ (Oxford) Postdoc, Visit Scholar on AI health care (e.g., ECG)
- ④ (Zenith Ai) Sr. Researcher on Omics AI (e.g., DNA, tRNA, protein, amino acid, ribosome, structure-sequence-function)

Technical novelty and principle

Example weighting is universal in deep learning

We define our interpretation of example weighting [6]:

Definition (Example Weighting). *In gradient-based optimisation, the derivative of an example can be interpreted as its effect on the update of a model. Therefore, a derivative's magnitude function equals to a weighting scheme.*

Accordingly, a change of the derivative magnitude function, is implicitly equivalent to, modifying an example weighting scheme.

Intuitive research motivations:

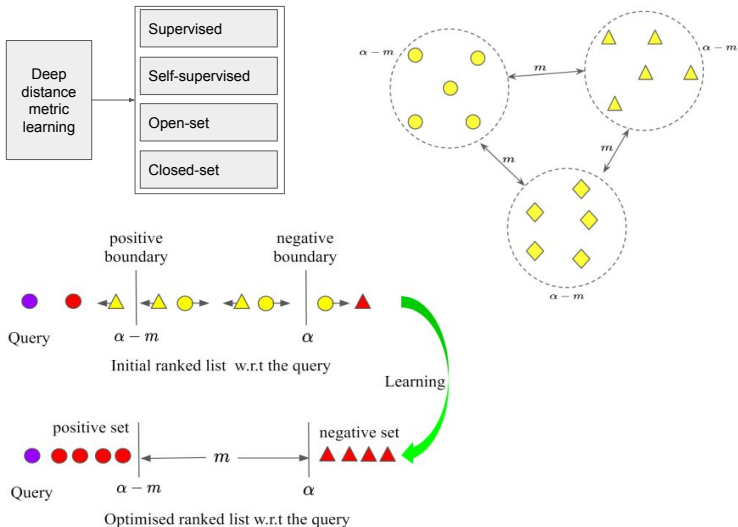
- (1) Not all training examples are created equal!
- (2) Sampling matters!

Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data**
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Structured learning on selected data

CVPR 2019 and TPAMI 2022 [8]



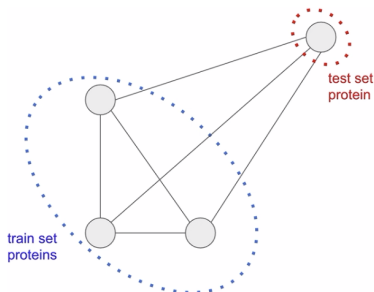
Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Learning objectives of robust deep learning

What is the meaning of robustness here?

- 1 To learn meaningful patterns on semantically clean data.
- 2 Without fitting errors/bias.
- 3 Generalisation to unseen data.



Build training and testing datasets properly [1].
How about the validation dataset?

Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Contribution: disclose the learning dynamics [7]

DNN has strong fitting capability, *but we find*:

Deep models easily fit random noise.

Deep networks learn simple semantic patterns before fitting noise.

Modern deep neural works tend to be over-confident.

(Ours: miscalibration under the noise) Deep neural networks become less confident of learning semantic patterns before fitting noise when the label noise rises.

Contribution: propose the ProSelfLC to promote confident and accurate learning [7]

1. To reward a low-entropy status other than penalise.
2. To promote model calibration.

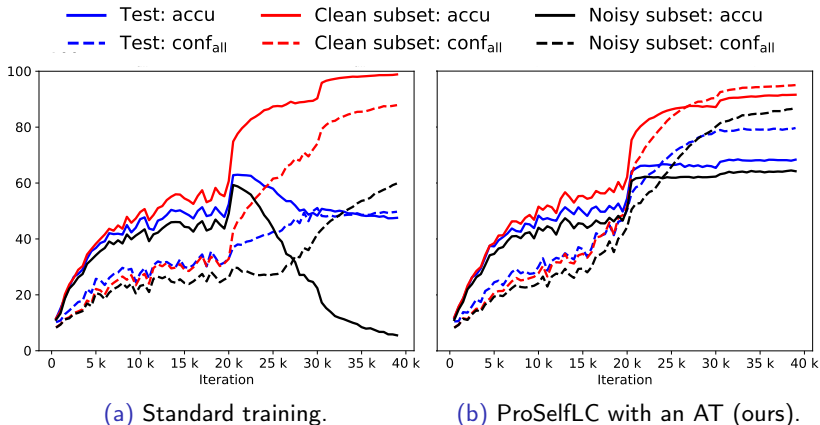


Figure: accu and conf_{all} when training ResNet18 on CIFAR-100. The noise rate is $r = 40\%$.

Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design**
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Cutting-edge OTS research: triplet design

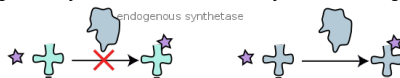
Orthogonal (aaRS, tRNA, ncAA)

With the conditions that

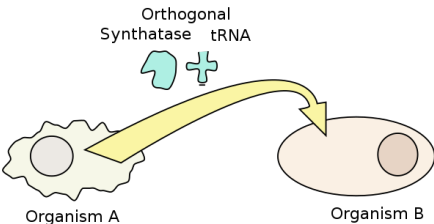
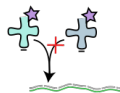
I
Orthogonal synthetase can aminoacylate only the orthogonal tRNA



II
Endogenous synthetases cannot aminoacylate the orthogonal tRNA

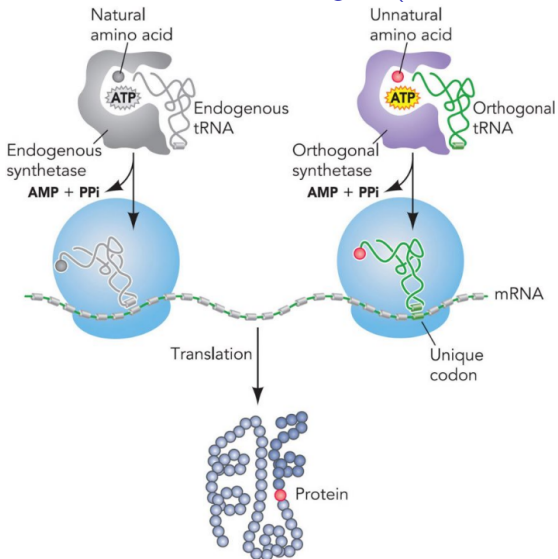


III
The orthogonal tRNA binds an unallocated codon



Cutting-edge OTS research: triplet design

Orthogonal (aaRS, tRNA, ncAA)



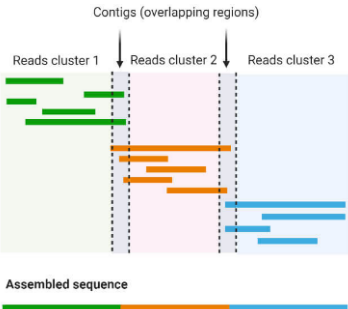
Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening**
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Bioinformatics

Sequence alignment and distance metrics

- 1 Gene synthesis: seq. optimisation + overlap-based assembly
- 2 NGS data analysis (e.g., variants calling + expression analysis)
 - Alignment: Bowtie2/minimap2 + Samtools
 - Reads counting
- 3 Alignment-based distance metrics
 - MMSeq2
 - Pfam domain database + HMMer-based distance



Domain organisation

Below is a listing of the unique domain organisations or architectures in which this domain is found. [More...](#)

There are 10142 sequences with the following architecture: **tRNA-synt_1c**, **Anticodon_2**

[W3ARE9_9FIRM](#) [Lachnospiraceae bacterium JC7] Glutamate--tRNA ligase [ECO:0000256]HAMAP-Rule:MF_00022) (482 residues)

tRNA-synt_1c Anticodon_2

[Show](#) all sequences with this architecture.

There are 5438 sequences with the following architecture: **tRNA-synt_1c**, **tRNA-synt_1c_C**

[YQKIT4_9PROT](#) [Methylophilaceae bacterium 11] Glutamine--tRNA ligase [ECO:0000256]HAMAP-Rule:MF_00126) (588 residues)

tRNA-synt_1c tRNA-synt_1c_C

[Show](#) all sequences with this architecture.

There are 3938 sequences with the following architecture: **tRNA-synt_1c**

[W9R803_9ROSA](#) [Morus notabilis] Glutamyl-tRNA synthetase [ECO:0000256]ARBA:ARBA00017458) (570 residues)

tRNA-synt_1c

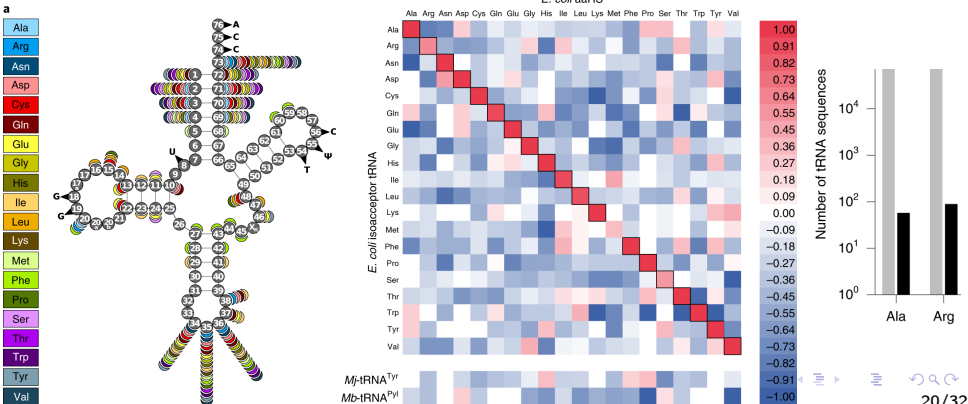
Bioinformatics for virtual screening

Revisiting Neil's question on sequences from diverse organisms

1 tRNA distance calculation [2]

- Secondary structure and the canonical numbering scheme.
- Identity elements, responsible for recognising cognate aaRS.

2 Annotating/validating via wetlab experiments: differential gene expression analysis.



Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling**
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Not-solved ideas: hierarchical objectives

StructureMotif-Loop-Site: Sequence-Domain-Site

Partial information

LEAFEKALKEM



????????????



LE????KA??EM



LEAF????KEM



Design task

Structure prediction

Sequence design

Loop design

Functional site design

Sequence

LEAFEKALKEM

Predicted Structure



Final Design



MCMC or gradient update

Loss function
- Hallucination
- Motif
- Problem-specific

Desired Motif



Partial Sequence

LEAF????KEM

Completed Sequence

LEAFEKALKEM

RF_{joint}
Neural network



Partial Structure

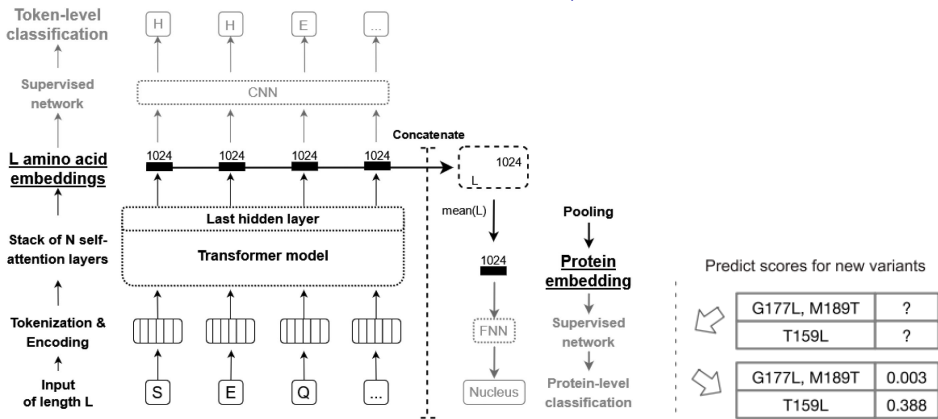


Completed Structure

[5] "Scaffolding protein functional sites using deep learning" Science (2022).

Mutational effect: single and co-mutation

Active AI for selective/evolutionary screening



[3] "ProtTrans: Toward Understanding the Language of Life Through Self-Supervised Learning." IEEE Transactions on Pattern Analysis and Machine Intelligence (2022).

[4] "Neural networks to learn protein sequence–function relationships from deep mutational scanning data." Proceedings of the National Academy of Sciences (2021).

[7] "ProSelfLC: Progressive Self Label Correction Towards A Low-Temperature Entropy State." **Ours** under peer review.

Public/Proprietary data collection, curation

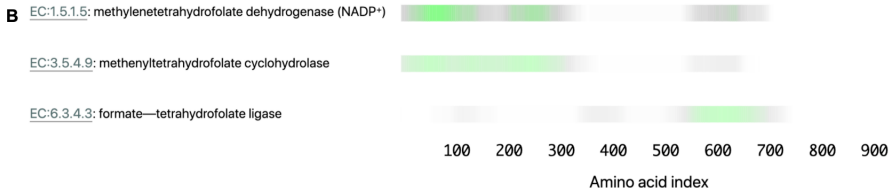
Interpretative ML: backbone + functional loops/regions/sites

Domain visualisation of the same EC function:



UniProt's annotations versus model's predictions:

A	Position(s)	Description	Graphical view	Length
	2 – 305	Methylenetetrahydrofolate dehydrogenase and cyclohydrolase		304
	306 – 935	Formyltetrahydrofolate synthetase		630



Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

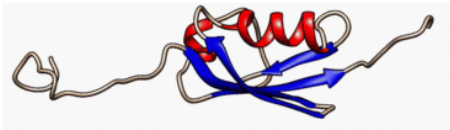
The key: collect and leverage the data

Data modelling pipelines based on alignment and similarity:

- Diversified datasets splitting: training, validation, testing
- A NN classifier as a baseline for the deep learning.

Data augmentation pipelines:

- Label propagation and transformation: A protein -> HMMScan -> domain annotations -> statistical association modelling -> EC/GO labels.
- Mutations in disordered regions have negligible effect, except for mutations to phenylalanine (P), tyrosine (T) and tryptophan (W), which promote order.



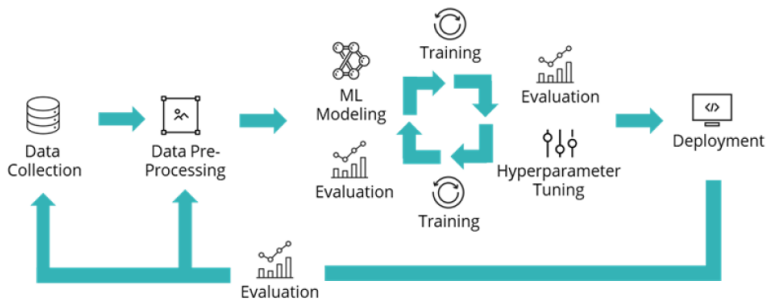
Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Industrial R&D for real-world problems

Methods = model-centric AI + data-centric AI

- Implement SOTA solutions & productise them.
- Build modularised AI toolboxes, end-to-end AI service pipelines.



Source: <https://dida.do/blog/data-centric-machine-learning>

Data collection, curation, and pre-processing take ~95% of the effort, and are harder to automate.

Outline

- 1 Self introduction
 - Overview
 - Work and experience
- 2 Technical novelty and principle: Example weighting is universal
- 3 Deep distance metric learning
 - Contribution: Structured learning on selected informative data
- 4 Robust deep learning (scope: transparent and interpretable ML)
 - Learning objectives
 - Contributions: an insightful finding and propose ProSelfLC
- 5 Omics AI
 - Cutting-edge OTS research: triplet design
 - Bioinformatics: alignment-based distance and virtual screening
 - Robust deep learning for protein modelling
 - The key: collecting, curating, and leveraging the data
- 6 Industrial R&D experience
 - Industrial R&D for real-world problems
 - R&D leading experience

Industrial AI Research

Processes and management

- 1 Suggest research **directions** and write **proposals** to the board (CEO and CSO).
- 2 Plan and Lead research, via **task breakdown** and an agreed completion **timeline** tracked by Jira and Confluence.
- 3 Collaborate, control quality, and maintain conventions using the **peer review** process, for both **code** and **documentation**.

Thanks for your attention.

Questions and discussions are very welcome.

Research topics and interests of Dr. (Amos) Xinshao Wang:

- Deep distance metric learning
- Robust deep learning
- Omics AI + Bioinformatics
- Active learning
- EDA + Data Visualisation

Google Scholar: [yOBhB7UAAAAJ](https://scholar.google.com/citations?hl=en&user=yOBhB7UAAAAJ)

Github: <https://github.com/XinshaoAmosWang>

Homepage+Blogs: <https://xinshaoamoswang.github.io/>

References

- [1] Bileschi, M. L., Belanger, D., Bryant, D. H., Sanderson, T., Carter, B., Sculley, D., Bateman, A., DePristo, M. A., and Colwell, L. J. Using deep learning to annotate the protein universe. *Nature Biotechnology*, pp. 1–6, 2022.
- [2] Cervettini, D., Tang, S., Fried, S. D., Willis, J. C., Funke, L. F., Colwell, L. J., and Chin, J. W. Rapid discovery and evolution of orthogonal aminoacyl-trna synthetase-trna pairs. *Nature biotechnology*, pp. 989–999, 2020.
- [3] Elnaggar, A., Heinzinger, M., Dallago, C., Rehawi, G., Yu, W., Jones, L., Gibbs, T., Feher, T., Angerer, C., Steinegger, M., et al. ProtTrans: Toward understanding the language of life through self-supervised learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 7112–7112.
- [4] Gelman, S., Fahlberg, S. A., Heinzelman, P., Romero, P. A., and Gitter, A. Neural networks to learn protein sequence–function relationships from deep mutational scanning data. *Proceedings of the National Academy of Sciences*, 2021.
- [5] Wang, J., Lianza, S., Juergens, D., Tischer, D., Watson, J. L., Castro, K. M., Ragotte, R., Saragovi, A., Milles, L. F., Baek, M., et al. Scaffolding protein functional sites using deep learning. *Science*, pp. 387–394, 2022.
- [6] Wang, X., Kodirov, E., Hua, Y., and Robertson, N. M. Derivative manipulation for general example weighting. *arXiv preprint arXiv:1905.11233*, 2019.
- [7] Wang, X., Hua, Y., Kodirov, E., Mukherjee, S. S., Clifton, D. A., and Robertson, N. M. Proselflc: Progressive self label correction towards a low-temperature entropy state. *arXiv preprint arXiv:2207.00118*, 2022.
- [8] Wang, X., Hua, Y., Kodirov, E., and Robertson, N. Ranked list loss for deep metric learning. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 5414–5429, 2022.